



INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification ⁷ : G10L		A2	(11) International Publication Number: WO 00/46787
			(43) International Publication Date: 10 August 2000 (10.08.00)
(21) International Application Number: PCT/US00/02808		(81) Designated States: AE, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CR, CU, CZ, DE, DK, DM, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, TZ, UA, UG, US, UZ, VN, YU, ZA, ZW, ARIPO patent (GH, GM, KE, LS, MW, SD, SL, SZ, TZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).	
(22) International Filing Date: 4 February 2000 (04.02.00)			
(30) Priority Data: 60/118,949 5 February 1999 (05.02.99) US			
(71) Applicant (for all designated States except US): CUSTOM SPEECH USA, INC. [US/US]; Suite B365, 3 North Court Street, Crown Point, IN 46307 (US).			
(72) Inventors; and			
(75) Inventors/Applicants (for US only): KAHN, Jonathan [US/US]; 1108 Cheyenne Drive, Crown Point, IN 46307 (US). QIN, Charles [-/US]; 23461 North Garden Lane, Lake Zurich, IL 60047 (US). FLYNN, Thomas, P. [US/US]; 562 Ridgelawn Road, Crown Point, IN 46307 (US). TIPPE, Robert, J. [US/US]; 3813 W. 214th Street, Matteson, IL 60443 (US).			
(74) Agents: SIGALE, Jordan, A. et al.; Sonnenschein, Nath & Rosenthal, 8000 Sears Tower, 233 S. Wacker Drive, Chicago, IL 60606-6404 (US).			

Published

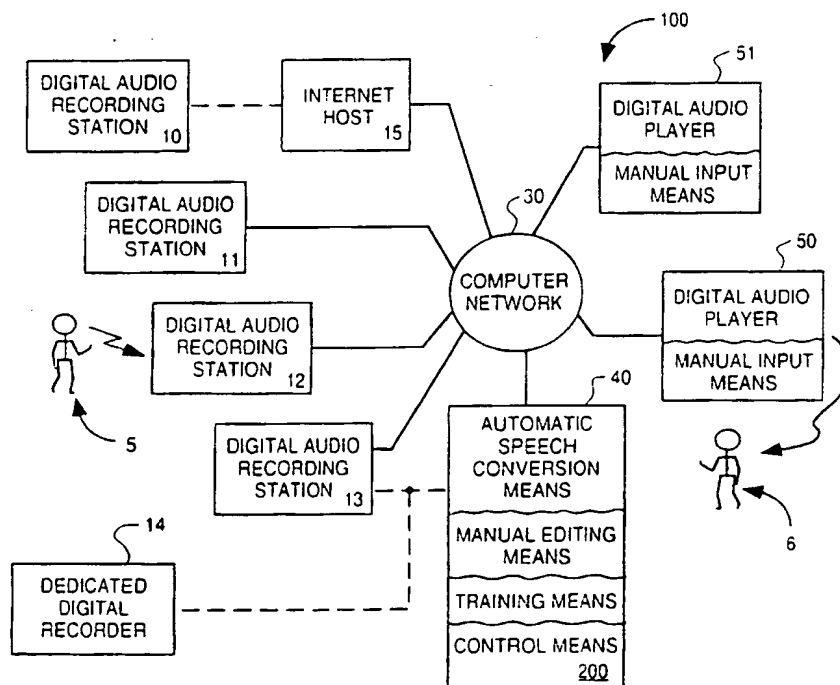
Without international search report and to be republished upon receipt of that report.

(54) Title: SYSTEM AND METHOD FOR AUTOMATING TRANSCRIPTION SERVICES

(57) Abstract

A system for substantially automating transcription services for multiple voice users including a manual transcription station, a speech recognition program and a routing program. The system establishes a profile for each of the voice users containing a training status which is selected from the group of enrollment, training, automated and stop automation. The system generates a uniquely identified voice dictation file from a current voice user and — based on the training status the system — routes the uniquely identified voice dictation file to a manual transcription station and the speech recognition program. A human transcriptionist creates transcribed files for each received voice dictation file. The speech recognition program automatically creates a written text for each received voice dictation file if the training status of the current user is training or automated. A verbatim file is manually established if the training status of the current user is enrollment or training and the speech recognition program is trained with an acoustic model for the current user using the verbatim file and the voice dictation file if the training status of the current user is enrollment or training.

The transcribed file is returned to the current user if the training status of the current user is enrollment or training or the written text is returned if the training status of the current user is automated. An apparatus and method is also disclosed for testing the skills of a transcriptionist. These apparatuses and methods may also be used to establish new base language model for mass distribution.



FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece	ML	Mali	TR	Turkey
BG	Bulgaria	HU	Hungary	MN	Mongolia	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MR	Mauritania	UA	Ukraine
BR	Brazil	IL	Israel	MW	Malawi	UG	Uganda
BY	Belarus	IS	Iceland	MX	Mexico	US	United States of America
CA	Canada	IT	Italy	NE	Niger	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NL	Netherlands	VN	Viet Nam
CG	Congo	KE	Kenya	NO	Norway	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NZ	New Zealand	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	PL	Poland		
CM	Cameroon	KR	Republic of Korea	PT	Portugal		
CN	China	KZ	Kazakstan	RO	Romania		
CU	Cuba	LC	Saint Lucia	RU	Russian Federation		
CZ	Czech Republic	LI	Liechtenstein	SD	Sudan		
DE	Germany	LK	Sri Lanka	SE	Sweden		
DK	Denmark	LR	Liberia	SG	Singapore		
EE	Estonia						

SYSTEM AND METHOD FOR AUTOMATING TRANSCRIPTION SERVICESBackground of the Invention

1. Field of the Invention

5 The present invention relates in general to computer speech recognition systems and, in particular, to a system and method for automating the text transcription of voice dictation by various end users.

2. Background Art

Speech recognition programs are well known in the art. While these programs are ultimately useful in automatically converting speech into text, many users are
10 dissuaded from using these programs because they require each user to spend a significant amount of time training the system. Usually this training begins by having each user read a series of pre-selected materials for approximately 20 minutes. Then, as the user continues to use the program, as words are improperly transcribed the user is expected to stop and train the program as to the intended word thus advancing the
15 ultimate accuracy of the acoustic model. Unfortunately, most professionals (doctors, dentists, veterinarians, lawyers) and business executive are unwilling to spend the time developing the necessary acoustic model to truly benefit from the automated transcription.

Accordingly, it is an object of the present invention to provide a system that
20 offers transparent training of the speech recognition program to the end-users.

There are systems for using computers for routing transcription from a group of end users. Most often these systems are used in large multi-user settings such as hospitals. In those systems, a voice user dictates into a general-purpose computer or other recording device and the resulting file is transferred automatically to a human
25 transcriptionist. The human transcriptionist transcribes the file, which is then returned to the original "author" for review. These systems have the perpetual overhead of employing a sufficient number of human transcriptionist to transcribe all of the dictation files.

PCT/US00/02808
Accordingly it is another object of the present invention to provide an automated means of translating speech into text where ever suitable so as to minimize the number of human transcriptionist necessary to transcribe audio files coming into the system.

5 It is an associated object to provide a simplified means for providing verbatim text files for training a user's acoustic model for the speech recognition portion of the system.

It is another associated object of the present invention to automate a preexisting speech recognition program toward further minimizing the number operators necessary to operate the system.

10 These and other objects will be apparent to those of ordinary skill in the art having the present drawings, specification and claims before them.

Summary of the Invention

The present invention comprises, in part, a system for substantially automating transcription services for one or more voice users. The system includes means for
15 creating a uniquely identified voice dictation file from a current user and an audio player used to audibly reproduce said uniquely identified voice dictation file. Both of these system elements can be implemented on the same or different general-purpose computers. Additionally, the voice dictation file creating means includes a system for assigning unique file handles to audio files and an audio recorder, and further comprise
20 means for operably connecting to a separate digital recording device and/or means for reading audio files from removable magnetic and other computer media.

Each of the general purpose computers implementing the system may be remotely located from the other computers but in operable connection to each other by way of a computer network, direct telephone connection, via email or other Internet
25 based transfer.

The system further includes means for manually inputting and creating a transcribed file based on humanly perceived contents of the uniquely identified voice dictation file. Thus, for certain voice dictation files, a human transcriptionist manually transcribes a textual version of the audio -- using a text editor or word processor -- based
30 on the output of the output of the audio player.

The system also includes means for automatically converting the voice dictation file into written text. The automatic speech converting means may be a preexisting speech recognition program, such as Dragon Systems' Naturally Speaking, IBM's Via Voice or Philips Corporation's Magic Speech. In such a case, the automatic speech
5 converting means includes means for automating responses to a series of interactive inquiries from the preexisting speech recognition program. In one embodiment, the system also includes means for manually selecting a specialized language model.

The system further includes means for manually editing the resulting written text to create a verbatim text of the voice dictation file. At the outset of a users use of the
10 system, this verbatim text will have to be created completely manually. However, after the automatic speech converting means has begun to sufficiently develop that user's acoustic model a more automated means can be used.

In a preferred embodiment, that manual editing means includes means for sequentially comparing a copy of the written text with the transcribed file resulting in a
15 sequential list of unmatched words culled from the copy of said written text. The manual editing means further includes means for incrementally searching for the current unmatched word contemporaneously within a first buffer associated with the speech recognition program containing the written text and a second buffer associated with the sequential list. Finally, the preferred manual editing means includes means for
20 correcting the current unmatched word in the second buffer, which includes means for displaying the current unmatched word in a manner substantially visually isolated from other text in the written text and means for playing a portion of the voice dictation recording from said first buffer associated with said current unmatched word. In one embodiment, the manual input means further includes means for alternatively viewing
25 the current unmatched word in context within the written text. For instance, the operator may wish to view the unmatched within the sentence in which it appears or perhaps with only is immediately adjacent words. Thus, the manner substantially visual isolation can be manually selected from the group containing word-by-word display, sentence-by-sentence display, and said current unmatched word display. The manual editing means
30 portion of the complete system may also be utilized as a separate apparatus.

The system may also include means for determining the skill of a human transcriptionist. In one approach, this accuracy determination can be made by

determining the ratio of the number of words in the sequential list of unmatched words to the number of words in the written text.

The system additionally includes means for training the automatic speech converting means to achieve higher accuracy for the current user. In particular, the training means uses the verbatim text created by the manual editing means and the voice dictation file. The training means may also comprise a preexisting training portion of the preexisting speech recognition program. Thus, the training means would also include means for automating responses to a series of interactive inquiries from the preexisting training portion of the speech recognition program. This functionality can be used, for instance, to establish a new language model (i.e. foreign language).

The system finally includes means for controlling the flow of the voice dictation file based upon the training status of the current user using the unique identification. The control means reads and modifies a user's training status such that it is an appropriate selection from the group of pre-enrollment, enrollment, training, automation and stop automation. During a user's pre-enrollment phase the control means further includes means for creating a user identification and acoustic model within the automatic speech converting means. During this phase, the control means routes the voice dictation file to the automatic speech converting means and the manual input means, routes the written text and the transcribed file to the manual editing means, routes the verbatim text to the training means and routes the transcribed file back to the current user as a finished text.

During the training phase, the control means routes (1) the voice dictation file to the automatic speech converting means and the manual input means, (2) routes the written text and the transcribed file to the manual editing means, (3) routes the verbatim text to the training means and (4) routes the transcribed file back to the current user as a finished text.

During the automation stage, the control means routes (1) the voice dictation file only to the automatic speech converting means and (2) the written text back to the current user as a finished text.

The present application also discloses a method for automating transcription services for one or more voice users in a system including a manual transcription station

and a speech recognition program. The method comprising the steps of: (1) establishing a profile for each of the voice users, the profile containing a training status; (2) creating a uniquely identified voice dictation file from a current voice user; (3) choosing the training status of the current voice user from the group of enrollment, training, automated and stop automation; (4) routing the voice dictation file to at least one of the manual transcription station and the speech recognition program based on the training status; (5) receiving the voice dictation file in at least one of the manual transcription station and the speech recognition program; (6) creating a transcribed file at the manual transcription station for each received voice dictation file; (7) automatically creating a written text with the speech recognition program for each received voice dictation file if the training status of the current user is training or automated; (8) manually establishing a verbatim file if the training status of the current user is enrollment or training; (9) training the speech recognition program with an acoustic model for the current user using the verbatim file and the voice dictation file if the training status of the current user is enrollment or training; (10) returning the transcribed file to the current user if the training status of the current user is enrollment or training; and (11) returning the written text to the current user if the training status of the current user is automated.

Brief Description of the Drawings

Fig. 1 of the drawings is a block diagram of one potential embodiment of the present system for substantially automating transcription services for one or more voice users;

Fig. 1b of the drawings is a block diagram of a general-purpose computer which may be used as a dictation station, a transcription station and the control means within the present system;

Fig. 2a of the drawings is a flow diagram of the main loop of the control means of the present system;

Fig. 2b of the drawings is a flow diagram of the enrollment stage portion of the control means of the present system;

Fig. 2c of the drawings is a flow diagram of the training stage portion of the control means of the present system;

WO 00/40787

Fig. 2d of the drawings is a flow diagram of the automation stage portion of the control means of the present system;

Fig. 3 of the drawings is a directory structure used by the control means in the present system;

5 Fig. 4 of the drawings is a block diagram of a portion of a preferred embodiment of the manual editing means; and

Fig. 5 of the drawings is an elevation view of the remainder of a preferred embodiment of the manual editing means.

Best Modes of Practicing the Invention

10 While the present invention may be embodied in many different forms, there is shown in the drawings and discussed herein a few specific embodiments with the understanding that the present disclosure is to be considered only as an exemplification of the principles of the invention and is not intended to limit the invention to the embodiments illustrated.

15 Fig. 1 of the drawings generally shows one potential embodiment of the present system for substantially automating transcription services for one or more voice users. The present system must include some means for receiving a voice dictation file from a current user. This voice dictation file receiving means can be a digital audio recorder, an analog audio recorder, or standard means for receiving computer files on magnetic media
20 or via a data connection.

As shown, in one embodiment, the system 100 includes multiple digital recording stations 10, 11, 12 and 13. Each digital recording station has at least a digital audio recorder and means for identifying the current voice user.

25 Preferably, each of these digital recording stations is implemented on a general-purpose computer (such as computer 20), although a specialized computer could be developed for this specific purpose. The general-purpose computer, though has the added advantage of being adaptable to varying uses in addition to operating within the present system 100. In general, the general-purpose computer should have, among other elements, a microprocessor (such as the Intel Corporation PENTIUM, Cyrix K6 or

Motorola 68000 series); volatile and non-volatile memory; one or more mass storage devices (i.e. HDD (not shown), floppy drive 21, and other removable media devices 22 such as a CD-ROM drive, DITTO, ZIP or JAZ drive (from Iomega Corporation) and the like); various user input devices, such as a mouse 23, a keyboard 24, or a microphone 25; and a video display system 26. In one embodiment, the general-purpose computer is controlled by the WINDOWS 9.x operating system. It is contemplated, however, that the present system would work equally well using a MACINTOSH computer or even another operating system such as a WINDOWS CE, UNIX or a JAVA based operating system, to name a few.

Regardless of the particular computer platform used, in an embodiment utilizing an analog audio input (via microphone 25) the general-purpose computer must include a sound-card (not shown). Of course, in an embodiment with a digital input no sound card would be necessary.

In the embodiment shown in Fig. 1, digital audio recording stations 10, 11, 12 and 13 are loaded and configured to run digital audio recording software on a PENTIUM-based computer system operating under WINDOWS 9.x. Such digital recording software is available as a utility in the WINDOWS 9.x operating system or from various third party vendor such as The Programmers' Consortium, Inc. of Oakton, Virginia (VOICEDOC), Syntrillium Corporation of Phoenix, Arizona (COOL EDIT) or Dragon Systems Corporation (Dragon Naturally Speaking Professional Edition). These various software programs produce a voice dictation file in the form of a "WAV" file. However, as would be known to those skilled in the art, other audio file formats, such as MP3 or DSS, could also be used to format the voice dictation file, without departing from the spirit of the present invention. In one embodiment where VOICEDOC software is used that software also automatically assigns a file handle to the WAV file, however, it would be known to those of ordinary skill in the art to save an audio file on a computer system using standard operating system file management methods.

Another means for receiving a voice dictation file is dedicated digital recorder 14, such as the Olympus Digital Voice Recorder D-1000 manufactured by the Olympus Corporation. Thus, if the current voice user is more comfortable with a more conventional type of dictation device, they can continue to use a dedicated digital recorder 14. In order to harvest the digital audio text file, upon completion of a

recording, dedicated digital recorder 14 would be operably connected to one of the digital audio recording stations, such as 13, toward downloading the digital audio file into that general-purpose computer. With this approach, for instance, no audio card would be required.

- 5 Another alternative for receiving the voice dictation file may consist of using one form or another of removable magnetic media containing a pre-recorded audio file. With this alternative an operator would input the removable magnetic media into one of the digital audio recording stations toward uploading the audio file into the system.

- 10 In some cases it may be necessary to pre-process the audio files to make them acceptable for processing by the speech recognition software. For instance, a DSS file format may have to be changed to a WAV file format, or the sampling rate of a digital audio file may have to be upsampled or downsampled. For instance, in use the Olympus Digital Voice Recorder with Dragon Naturally Speaking, Olympus' 8MHz rate needs to be upsampled to 11 MHz. Software to accomplish such pre-processing is available from
15 a variety of sources including Syntrillium Corporation and Olympus Corporation.

- The other aspect of the digital audio recording stations is some means for identifying the current voice user. The identifying means may include keyboard 24 upon which the user (or a separate operator) can input the current user's unique identification code. Of course, the user identification can be input using a myriad of computer input
20 devices such as pointing devices (e.g. mouse 23), a touch screen (not shown), a light pen (not shown), bar-code reader (not shown) or audio cues via microphone 25, to name a few.

- In the case of a first time user the identifying means may also assign that user an identification number after receiving potentially identifying information from that user, including: (1) name; (2) address; (3) occupation; (4) vocal dialect or accent; etc. As
25 discussed in association with the control means, based upon this input information, a voice user profile and a sub-directory within the control means are established. Thus, regardless of the particular identification means used, a user identification must be established for each voice user and subsequently provided with a corresponding digital
30 audio file for each use such that the control means can appropriately route and the system ultimately transcribe the audio.

In one embodiment of the present invention, the identifying means may also seek the manual selection of a specialty vocabulary. It is contemplated that the specialty vocabulary sets may be general for various users such as medical (i.e. Radiology, Orthopedic Surgery, Gynecology) and legal (i.e. corporate, patent, litigation) or highly specific such that within each specialty the vocabulary parameters could be further limited based on the particular circumstances of a particular dictation file. For instance, if the current voice user is a Radiologist dictating the reading of a abdominal CAT scan the nomenclature is highly specialized and different from the nomenclature for a renal ultrasound. By narrowly segmenting each selectable vocabulary set an increase in the accuracy of the automatic speech converter is likely.

As shown in Fig. 1, the digital audio recording stations may be operably connected to system 100 as part of computer network 30 or, alternatively, they may be operably connected to the system via internet host 15. As shown in Fig. 1b, the general-purpose computer can be connected to both network jack 27 and telephone jack. With the use of an internet host, connection may be accomplished by e-mailing the audio file via the Internet. Another method for completing such connection is by way of direct modem connection via remote control software, such as PC ANYWHERE, which is available from Symantec Corporation of Cupertino, California. It is also possible, if the IP address of digital audio recording station 10 or internet host 15 is known, to transfer the audio file using basic file transfer protocol. Thus, as can be seen from the foregoing, the present system allows great flexibility for voice users to provide audio input into the system.

Control means 200 controls the flow of voice dictation file based upon the training status of the current voice user. As shown in Figs. 2a, 2b, 2c, 2d, control means 200 comprises a software program operating on general purpose computer 40. In particular, the program is initialized in step 201 where variable are set, buffers cleared and the particular configuration for this particular installation of the control means is loaded. Control means continually monitors a target directory (such as "current" (shown in Fig. 3)) to determine whether a new file has been moved into the target, step 202. Once a new file is found (such as "6723.id" (shown in Fig. 3)), a determination is made as to whether or not the current user 5 (shown in Fig. 1) is a new user, step 203.

For each new user (as indicated by the existence of a "pro" file in the "current" subdirectory), a new subdirectory is established, step 204 (such as the "usern" subdirectory (shown in Fig. 3)). This subdirectory is used to store all of the audio files ("xxxx.wav"), written text ("xxxx.wrt"), verbatim text ("xxxx.vb"), transcription text ("xxxx.txt") and user profile ("usern.pro") for that particular user. Each particular job is assigned a unique number "xxxx" such that all of the files associated with a job can be associated by that number. With this directory structure, the number of users is practically limited only by storage space within general-purpose computer 40.

Now that the user subdirectory has been established, the user profile is moved to the subdirectory, step 205. The contents of this user profile may vary between systems. The contents of one potential user profile is shown in Fig. 3 as containing: the user name, address, occupation and training status. Aside from the training status variable, which is necessary, the other data is useful in routing and transcribing the audio files.

The control means, having selected one set of files by the handle, determines the identity of the current user by comparing the ".id" file with its "user.tbl," step 206. Now that the user is known the user profile may be parsed from that user's subdirectory and the current training status determined, step 207. Steps 208-211 are the triage of the current training status is one of: enrollment, training, automate, and stop automation.

Enrollment is the first stage in automating transcription services. As shown in Fig. 2b, the audio file is sent to transcription, step 301. In particular, the "xxxx.wav" file is transferred to transcriptionist stations 50 and 51. In a preferred embodiment, both stations are general-purpose computers, which run both an audio player and manual input means. The audio player is likely to be a digital audio player, although it is possible that an analog audio file could be transferred to the stations. Various audio players are commonly available including a utility in the WINDOWS 9.x operating system and various other third parties such from The Programmers' Consortium, Inc. of Oakton, Virginia (VOICESCRIBE). Regardless of the audio player used to play the audio file, manual input means is running on the computer at the same time. This manual input means may comprise any of text editor or word processor (such as MS WORD, WordPerfect, AmiPro or Word Pad) in combination with a keyboard, mouse, or other user-interface device. In one embodiment of the present invention, this manual input means may, itself, also be speech recognition software, such as Naturally Speaking from

Dragon Systems of Newton, Massachusetts, Via Voice from IBM Corporation of Armonk, New York, or Speech Magic from Philips Corporation of Atlanta, Georgia.

Human transcriptionist 6 listens to the audio file created by current user 5 and as is known, manually inputs the perceived contents of that recorded text, thus establishing the transcribed file, step 302. Being human, human transcriptionist 6 is likely to impose experience, education and biases on the text and thus not input a verbatim transcript of the audio file. Upon completion of the human transcription, the human transcriptionist 6 saves the file and indicates that it is ready for transfer to the current users subdirectory as "xxxx.txt", step 303.

10 Inasmuch as this current user is only at the enrollment stage, a human operator will have to listen to the audio file and manually compare it to the transcribed file and create a verbatim file, step 304. That verbatim file "xxxx.vb" is also transferred to the current user's subdirectory, step 305. Now that verbatim text is available, control means 200 starts the automatic speech conversion means, step 306. This automatic speech
15 conversion means may be a preexisting program, such as Dragon System's Naturally Speaking, IBM's Via Voice or Philips' Speech Magic, to name a few. Alternatively, it could be a unique program that is designed to specifically perform automated speech recognition.

In a preferred embodiment, Dragon Systems' Naturally Speaking has been used
20 by running an executable simultaneously with Naturally Speaking that feeds phantom keystrokes and mousing operations through the WIN32API, such that Naturally Speaking believes that it is interacting with a human being, when in fact it is being controlled by control means 200. Such techniques are well known in the computer software testing art and, thus, will not be discussed in detail. It should suffice to say that
25 by watching the application flow of any speech recognition program, an executable to mimic the interactive manual steps can be created.

If the current user is a new user, the speech recognition program will need to establish the new user, step 307. Control means provides the necessary information from the user profile found in the current user's subdirectory. All speech recognition require
30 significant training to establish an acoustic model of a particular user. In the case of Dragon, initially the program seeks approximately 20 minutes of audio usually obtained by the user reading a canned text provided by Dragon Systems. There is also

functionality built into Dragon that allows "mobile training." Using this feature, the verbatim file and audio file are fed into the speech recognition program to beginning training the acoustic model for that user, step 308. Regardless of the length of that audio file, control means 200 closes the speech recognition program at the completion of the file, step 309.

As the enrollment step is too soon to use the automatically created text, a copy of the transcribed file is sent to the current user using the address information contained in the user profile, step 310. This address can be a street address or an e-mail address. Following that transmission, the program returns to the main loop on Fig. 2a.

After a certain number of minutes of training have been conducted for a particular user, that user's training status may be changed from enrollment to training. The border for this change is subjective, but perhaps a good rule of thumb is once Dragon appears to be creating written text with 80% accuracy or more, the switch between states can be made. Thus, for such a user the next transcription event will prompt control means 200 into the training state. As shown in Fig. 2c, steps 401-403 are the same human transcription steps as steps 301-303 in the enrollment phase. Once the transcribed file is established, control means 200 starts the automatic speech conversion means (or speech recognition program) and selects the current user, step 404. The audio file is fed into the speech recognition program and a written text is established within the program buffer, step 405. In the case of Dragon, this buffer is given the same file handle on very instance of the program. Thus, that buffer can be easily copied using standard operating system commands and manual editing can begin, step 406.

In one particular embodiment utilizing the VOICEWARE system from The Programmers' Consortium, Inc. of Oakton, Virginia, the user inputs audio into the VOICEWARE system's VOICEDOC program, thus, creating a ".wav" file. In addition, before releasing this ".wav" file to the VOICEWARE server, the user selects a "transcriptionist." This "transcriptionist" may be a particular human transcriptionist or may be the "computerized transcriptionist." If the user selects a "computerized transcriptionist" they may also select whether that transcription is handled locally or remotely. This file is assigned a job number by the VOICEWARE server, which routes the job to the VOICESCRIBE portion of the system. Normally, VOICESCRIBE is used by the human transcriptionist to receive and playback the job's audio (".wav") file. In

addition, the audio file is grabbed by the automatic speech conversion means. In this VOICEWARE system embodiment, by placing VOICESCRIBE in "auto mode" new jobs (i.e. an audio file newly created by VOICEDOC) are automatically downloaded from the VOICEWARE server and a VOICESCRIBE window having a window title
5 formed by the job number of the current ".wav" file. An executable file, running in the background "sees" the VOICESCRIBE window open and using the WIN32API determines the job number from the VOICESCRIBE window title. The executable file then launches the automatic speech conversion means. In Dragon System's Naturally Speaking, for instance, there is a built in function for performing speech recognition on a
10 preexisting ".wav" file. The executable program feeds phantom keystrokes to Naturally Speaking to open the ".wav" file from the "current" directory (see Fig. 3) having the job number of the current job.

In this embodiment, after Naturally Speaking has completed automatically transcribing the contexts of the ".wav" file, the executable file resumes operation by
15 selecting all of the text in the open Naturally Speaking window and copying it to the WINDOWS 9.x operating system clipboard. Then, using the clipboard utility, save the clipboard as a text file using the current job number with a "dmt" suffix. The executable file then "clicks" the "complete" button in VOICESCRIBE to return the "dmt" file to the VOICEWARE server. As would be understood by those of ordinary skill in the art, the
20 foregoing procedure can be done utilizing other digital recording software and other automatic speech conversion means. Additionally, functionality analogous to the WINDOWS clipboard exists in other operating systems. It is also possible to require human intervention to activate or prompt one or more of the foregoing steps. Further, although, the various programs executing various steps of this could be running on a
25 number of interconnected computers (via a LAN, WAN, internet connectivity, email and the like), it is also contemplated that all of the necessary software can be running on a single computer.

Another alternative approach is also contemplated wherein the user dictates directly into the automatic speech conversion means and the VOICEWARE server picks
30 up a copy in the reverse direction. This approach works as follows; without actually recording any voice, the user clicks on the "complete" button in VOICEDOC, thus, creating an empty ".wav" file. This empty file is nevertheless assigned a unique job

number by the VOICEWARE server. The user (or an executable file running in the

background) then launches the automatic speech conversion means and the user dictates directly into that program, in the same manner previously used in association with such automatic speech conversion means. Upon completion of the dictation, the user presses
5 a button labeled "return" (generated by a background executable file), which executable then commences a macro that gets the current job number from VOICEWARE (in the manner describe above), selects all of the text in the document and copies it to the clipboard. The clipboard is then saved to the file "<jobnumber>.dmt," as discussed above. The executable then "clicks" the "complete" button (via the WIN32API) in
10 VOICESCRIBE, which effectively returns the automatically transcribed text file back to the VOICEWARE server, which, in turn, returns the completed transcription to the VOICESCRIBE user. Notably, although, the various programs executing various steps of this could be running on a number of interconnected computers (via a LAN, WAN, internet connectivity, email and the like), it is also contemplated that all of the necessary
15 software can be running on a single computer. As would be understood by those of ordinary skill in the art, the foregoing procedure can be done utilizing other digital recording software and other automatic speech conversion means. Additionally, functionality analogous to the WINDOWS clipboard exists in other operating systems. It is also possible to require human intervention to activate or prompt one or more of the
20 foregoing steps.

Manual editing is not an easy task. Human beings are prone to errors. Thus, the present invention also includes means for improving on that task. As shown in Fig. 4, the transcribed file ("3333.txt") and the copy of the written text ("3333.wrt") are sequentially compared word by word 406a toward establishing sequential list of
25 unmatched words 406b that are culled from the copy of the written text. This list has a beginning and an end and pointer 406c to the current unmatched word. Underlying the sequential list is another list of objects which contains the original unmatched words, as well as the words immediately before and after that unmatched word, the starting location in memory of each unmatched word in the sequential list of unmatched words
30 406b and the length of the unmatched word.

As shown in Fig. 5, the unmatched word pointed at by pointer 406c from list 406b is displayed in substantial visual isolation from the other text in the copy of the

written text on a standard computer monitor 500 in an active window 501. As shown in Fig. 5, the context of the unmatched word can be selected by the operator to be shown within the sentence it resides, word by word or in phrase context, by clicking on buttons 514, 515, and 516, respectively.

5 Associated with active window 501 is background window 502, which contains the copy of the written text file. As shown in background window 502, a incremental search has located (see pointer 503) the next occurrence of the current unmatched word "cash." Contemporaneously therewith, within window 505 containing the buffer from the speech recognition program, the same incremental search has located (see pointer
10 506) the next occurrence of the current unmatched word. A human user will likely only being viewing active window 501 activate the audio replay from the speech recognition program by clicking on "play" button 510, which plays the audio synchronized to the text at pointer 506. Based on that snippet of speech, which can be played over and over by clicking on the play button, the human user can manually input the correction to the
15 current unmatched word via keyboard, mousing actions, or possibly even audible cues to another speech recognition program running within this window.

In the present example, even if the choice of isolated context offered by buttons 514, 515 and 516, it may still be difficult to determine the correct verbatim word out-of-
context, accordingly there is a switch window button 513 that will move background
20 window 502 to the foreground with visible pointer 503 indicating the current location within the copy of the written text. The user can then return to the active window and input the correct word, "trash." This change will only effect the copy of the written text displayed in background window 502.

When the operator is ready for the next unmatched word, the operator clicks on
25 the advance button 511, which advances pointer 406c down the list of unmatched words and activates the incremental search in both window 502 and 505. This unmatched word is now displayed in isolation and the operator can play the synchronized speech from the speech recognition program and correct this word as well. If at any point in the operation, the operator would like to return to a previous unmatched word, the operator
30 clicks on the reverse button 512, which moves pointer 406c back a word in the list and causes a backward incremental search to occur. This is accomplished by using the underlying list of objects which contains the original unmatched words. This list is

traversed in object by object fashion, but alternatively each of the records could be padded such that each item has the same word size to assist in bi-directional traversing of the list. As the unmatched words in this underlying list are read only it is possible to return to the original unmatched word such that the operator can determine if a different
 5 correction should have been made.

Ultimately, the copy of the written text is finally corrected resulting in a verbatim copy, which is saved to the user's subdirectory. The verbatim file is also passed to the speech recognition program for training, step 407. The new (and improved) acoustic model is saved, step 408, and the speech recognition program is closed, step 409. As the
 10 system is still in training, the transcribed file is returned to the user, as in step 310 from the enrollment phase.

As shown in Fig. 4, the system may also include means for determining the accuracy rate from the output of the sequential comparing means. Specifically, by counting the number of words in the written text and the number of words in list 406b
 15 the ratio of words in said sequential list to words in said written text can be determined, thus providing an accuracy percentage. As before, it is a matter of choice when to advance users from one stage to another. Once that goal is reached, the user's profile is changed to the next stage, step 211.

One potential enhancement or derivative functionality is provided by the
 20 determination of the accuracy percentage. In one embodiment, this percentage could be used to evaluate a human transcriptionist's skills. In particular, by using either a known verbatim file or a well-established user, the associated ".wav" file would be played for the human transcriptionist and the foregoing comparison would be performed on the transcribed text versus the verbatim file created by the foregoing process. In this
 25 manner, additional functionality can be provided by the present system.

As understood, currently, manufacturers of speech recognition programs use recording of foreign languages, dictions, etc. with manually established verbatim files to program speech models. It should be readily apparent that the foregoing manner of establishing verbatim text could be used in the initial development of these speech files
 30 simplifying this process greatly.

Once the user has reached the automation stage, the greatest benefits of the present system can be achieved. The speech recognition software is started, step 600, and the current user selected, step 601. If desired, a particularized vocabulary may be selected, step 602. Then automatic conversion of the digital audio file recorded by the current user may commence, step 603. When completed, the written file is transmitted to the user based on the information contained in the user profile, step 604 and the program is returned to the main loop.

Unfortunately, there may be instances where the voice users cannot use automated transcription for a period of time (during an illness, after dental work, etc.) because their acoustic model has been temporarily (or even permanently) altered. In that case, the system administrator may set the training status variable to a stop automation state in which steps 301, 302, 303, 305 and 310 (see Fig. 2b) are the only steps performed.

The foregoing description and drawings merely explain and illustrate the invention and the invention is not limited thereto. Those of the skill in the art who have the disclosure before them will be able to make modifications and variations therein without departing from the scope of the present invention. For instance, it is possible to implement all of the elements of the present system on a single general-purpose computer by essentially time sharing the machine between the voice user, transcriptionist and the speech recognition program. The resulting cost saving makes this system accessible to more types of office situations not simply large medical clinics, hospital, law firms or other large entities.

WHAT IS CLAIMED IS:

1. A system for substantially automating transcription services for one or more voice users, said system comprising:
 - 5 - means for creating a uniquely identified voice dictation file from a current user, said current user being one of said one or more voice users;
 - an audio player used to audibly reproduce said uniquely identified voice dictation file;
 - means for manually inputting and creating a transcribed file based on
10 humanly perceived contents of said uniquely identified voice dictation file;
 - means for automatically converting said uniquely identified voice dictation file into written text;
 - means for manually editing a copy of said written text to create a verbatim
15 text of said uniquely identified voice dictation file;
 - means for training said automatic speech converting means to achieve higher accuracy with said uniquely identified voice dictation file of current user; and
 - means for controlling the flow of said uniquely identified voice dictation
20 file based upon a training status of said current user, whereby said controlling means sends said uniquely identified voice dictation file to at least one of said manual input means and said automatic speech converting means.

2. The invention according to Claim 1 further comprising means for transferring said written text into a written text file.
3. The invention according to Claim 1 wherein said written text is at least temporarily synchronized to said uniquely identified voice dictation file, said manual editing means comprises:
 - means for sequentially comparing a copy of said written text with said transcribed file resulting in a sequential list of unmatched words culled from said copy of said written text, said sequential list having a beginning, an end and a current unmatched word, said current unmatched word being successively advanced from said beginning to said end;
 - means for incrementally searching for said current unmatched word contemporaneously within a first buffer associated with the speech recognition program containing said written text and a second buffer associated with said sequential list; and
 - means for correcting said current unmatched word in said second buffer, said correcting means including means for displaying said current unmatched word in a manner substantially visually isolated from other text in said copy of said written text and means for playing a portion of said synchronized voice dictation recording from said first buffer associated with said current unmatched word.
4. The invention according to Claim 3 wherein said correcting means further includes means for alternatively viewing said current unmatched word in context within said copy of said written text.
5. The invention according to Claim 3 further including means for determining an accuracy rate for said current user.
6. The invention according to Claim 5 wherein said verbatim file is a known accurate file, invention further includes means for determining skill of a human operator based on said accuracy rate.

7. The invention according to Claim 4 wherein said sequential list and said written text each have a respective number of words, said accuracy rate determining means determines the ratio of said number of words in said sequential list to said number of words in said written text.
- 5 8. The invention according to Claim 1 wherein said automatic speech converting means comprises a preexisting speech recognition program intended for human interactive use, said automatic speech converting means includes means for automating responses to a series of interactive inquiries from said preexisting speech recognition program.
- 10 9. The invention according to Claim 8 wherein said training means comprises a preexisting training portion of said preexisting speech recognition program intended for human interactive use, said training means includes means for automating responses to a series of interactive inquiries from said preexisting training portion of said preexisting speech recognition program.
- 15 10. The invention according to Claim 1 wherein said training means comprises a preexisting training program intended for human interactive use, said training means includes means for automating responses to a series of interactive inquiries from said preexisting training program.
- 20 11. The invention according to Claim 1 wherein said control means reads and modifies a user profile associated with said current user, each of said user profiles including said training status of said current user.
12. The invention according to Claim 11 wherein said training status is selected from the group of pre-enrollment, enrollment, training, automation and stop automation.
- 25 13. The invention according to Claim 12 when said training status is pre-enrollment said control means further includes means for creating a user identification and acoustic model within said automatic speech converting means.
14. The invention according to Claim 12 when said training status is enrollment said control means routes said voice dictation file to said automatic speech converting

means and said manual input means, routes said written text and said transcribed file to said manual editing means, routes said verbatim text to said training means and routes said transcribed file back to said current user as a finished text.

15. The invention according to Claim 12 when said training status is training said control means routes said voice dictation file to said automatic speech converting means and said manual input means, routes said written text and said transcribed file to said manual editing means, routes said verbatim text to said training means and routes said transcribed file back to said current user as a finished text.
16. The invention according to Claim 12 when said training status is automation said control means routes said voice dictation file only to said automatic speech converting means and routes said written text back to said current user as a finished text.
17. An apparatus for substantially simplifying the production of a foreign language speech model for said speech recognition program wherein said foreign language provides a sufficient set of words to teach the voice dictation recording based upon a transcribed file produced by a human transcriptionist and a written text produced by a speech recognition program, wherein said written text is at least temporarily synchronized to said voice dictation recording, said apparatus comprising:
- means for sequentially comparing a copy of said written text with said transcribed file resulting in a sequential list of unmatched words culled from said copy of said written text, said sequential list having a beginning, an end and a current unmatched word, said current unmatched word being successively advanced from said beginning to said end;
 - means for incrementally searching for said current unmatched word contemporaneously within a first buffer associated with the speech recognition program containing said written text and a second buffer associated with said sequential list; and
 - means for correcting said current unmatched word in said second buffer, said correcting means including means for displaying said current

unmatched word in a manner substantially visually isolated from other text in said copy of said written text and means for playing a portion of said synchronized voice dictation recording from said first buffer associated with said current unmatched word.

- 5 18. The invention according to Claim 17 wherein said correcting means further includes means for alternatively viewing said current unmatched word in context within said copy of said written text.
- 10 19. The invention according to Claim 18 wherein said manner substantially visually isolated from other text can be manually selected from the group containing word-by-word display, sentence-by-sentence display, and said current unmatched word display.
20. A method for automating transcription services for one or more voice users in a system including a manual transcription station and a speech recognition program, said method comprising the steps of:
- 15 - establishing a profile for each of the voice users, the profile containing a training status;
- 20 - creating a uniquely identified voice dictation file for a current voice user;
- 20 - choosing the training status of the current voice user from the group of enrollment, training, automated and stop automation;
- 20 - routing the uniquely identified voice dictation file to at least one of the manual transcription station and the speech recognition program based on the training status;
- 25 - receiving the uniquely identified voice dictation file in at least one of the manual transcription station and the speech recognition program;
- 25 - creating a transcribed file at the manual transcription station for each received uniquely identified voice dictation file;

- automatically creating a written text with the speech recognition program for each received uniquely identified voice dictation file if the training status of the current user is training or automated;
- 5 - manually establishing a verbatim file if the training status of the current user is enrollment or training;
- training the speech recognition program with an acoustic model for the current user using the verbatim file and the uniquely identified voice dictation file if the training status of the current user is enrollment or training;
- 10 - returning the transcribed file to the current user if the training status of the current user is enrollment or training; and
- returning the written text to the current user if the training status of the current user is automated.

21. The invention according to Claim 30 wherein said step of manually establishing a
15 verbatim file includes the sub-steps of:

- assisting an operator to establish the verbatim file if the training status of the current user is training by:
 - sequentially comparing a copy of the written text with the
20 transcribed file resulting in a sequential list of unmatched words culled from the copy of the written text, the sequential list having a beginning, an end and a current unmatched word, the current unmatched word being successively advanced from the beginning to the end;
 - incrementally searching for the current unmatched word
25 contemporaneously within a first buffer associated with the speech recognition program containing the written text and a second buffer associated with the sequential list; and

5

- displaying the current unmatched word in a manner substantially visually isolated from other text in the copy of the written text and playing a portion of the synchronized voice dictation recording from the first buffer associated with the current unmatched word; and
- correcting the current unmatched word to be a verbatim representation of the portion of the synchronized voice dictation recording.

10

22. A method for testing the skills of a human transcriptionist using a known accurate written text created by a speech recognition program and a transcribed file created by the human transcriptionist, the method comprising:

15

- sequentially comparing a copy of the written text with the transcribed file resulting in a sequential list of unmatched words culled from the copy of the written text, the sequential list having a beginning, an end and a current unmatched word, the current unmatched word being successively advanced from the beginning to the end;

20

- incrementally searching for the current unmatched word contemporaneously within a first buffer associated with the speech recognition program containing the written text and a second buffer associated with the sequential list; and

25

- displaying the current unmatched word in a manner substantially visually isolated from other text in the copy of the written text and playing a portion of the synchronized voice dictation recording from the first buffer associated with the current unmatched word; and
- calculating the accuracy rate of the human transcriptionist.

1/7

FIG. 1

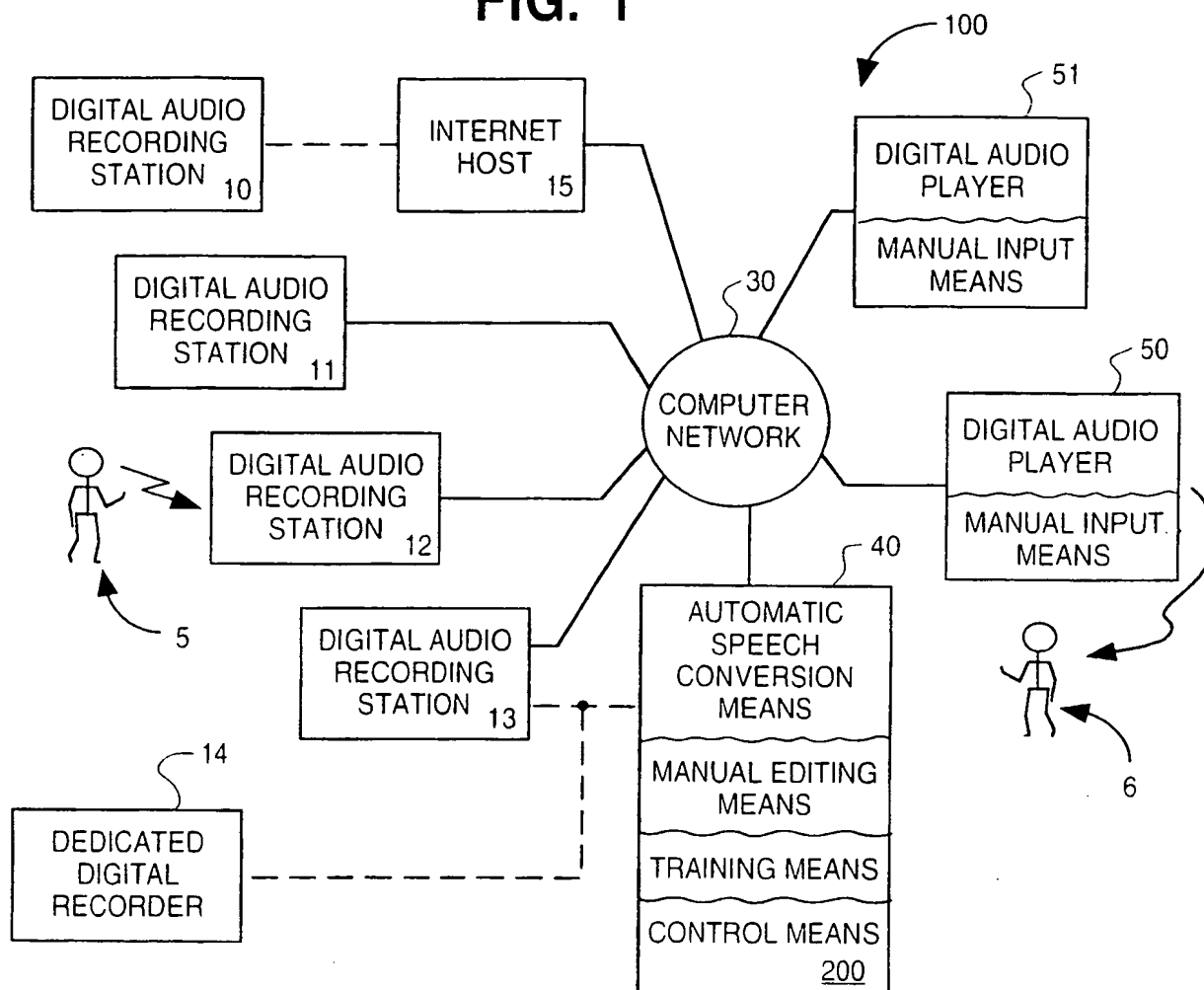


FIG. 1b

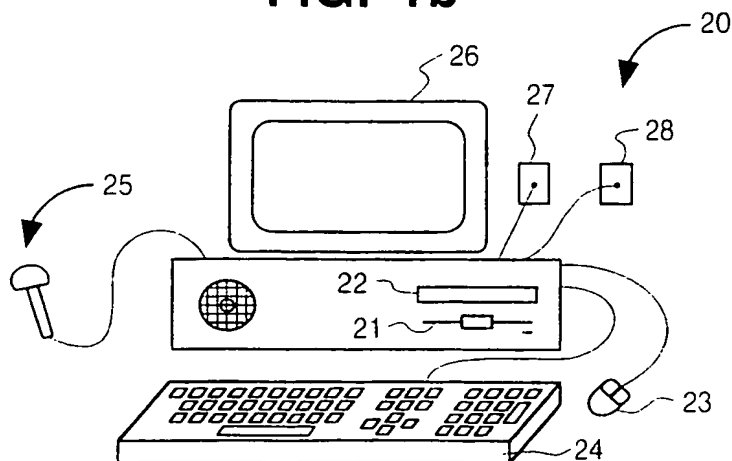


FIG. 2a

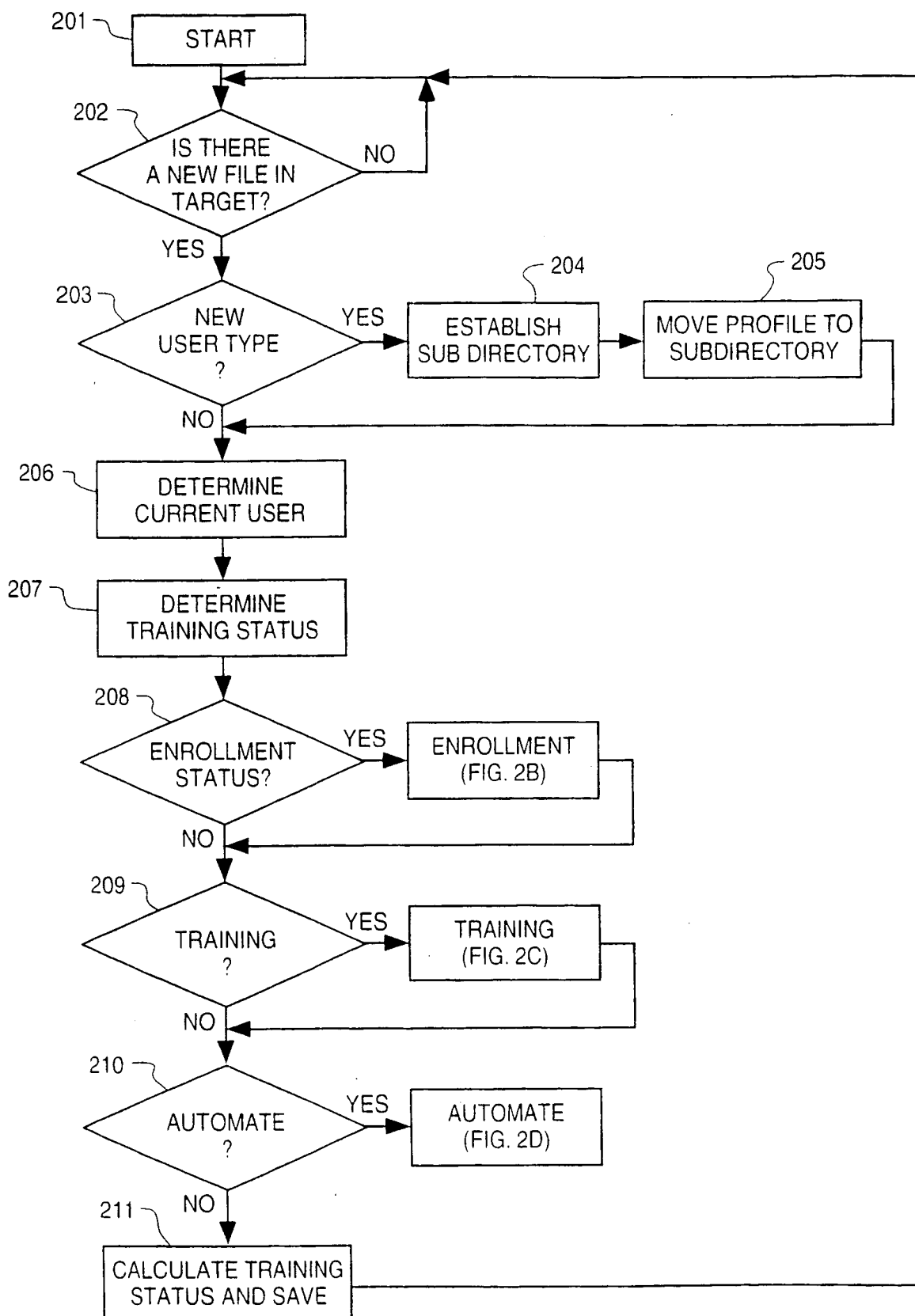


FIG. 2b

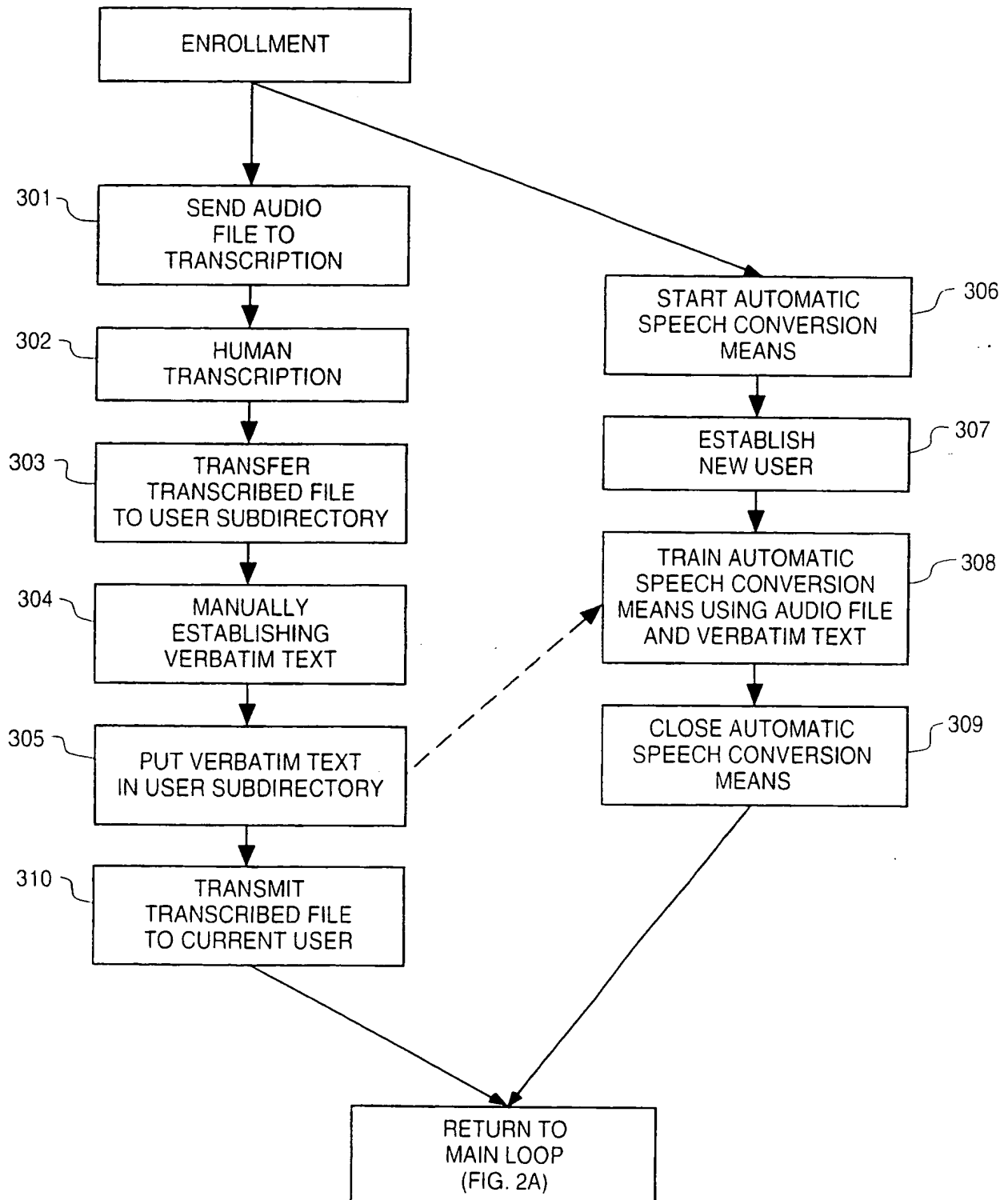


FIG. 2c

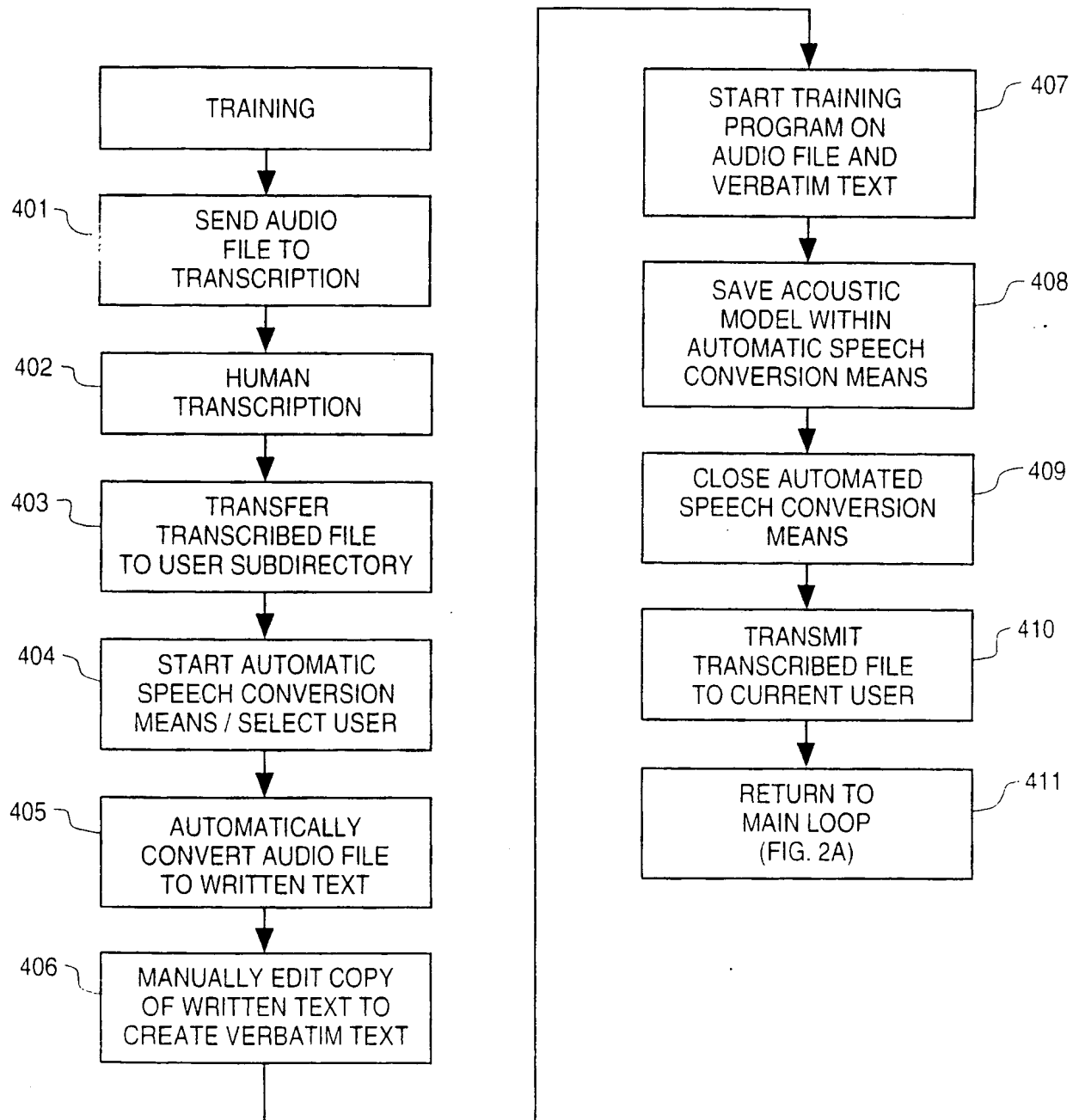
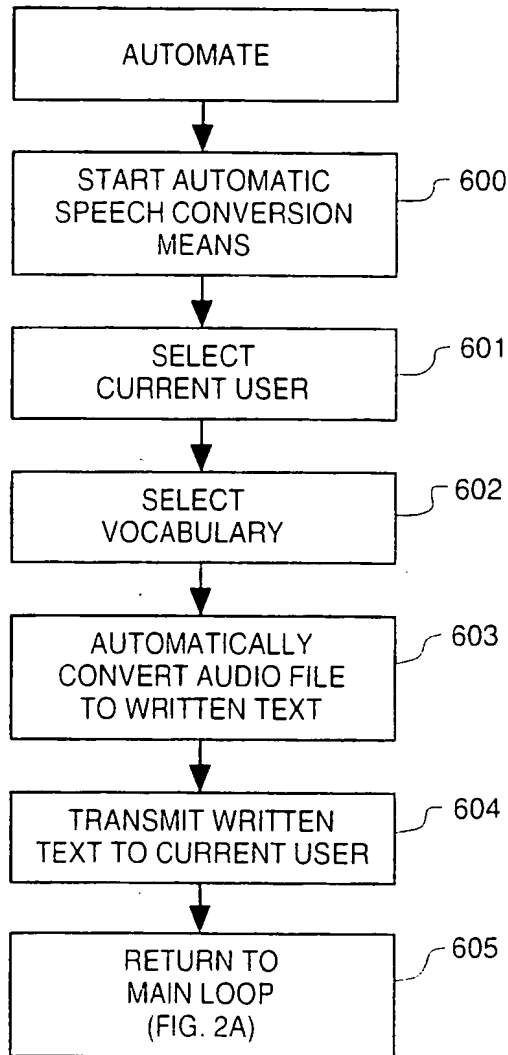


FIG. 2d

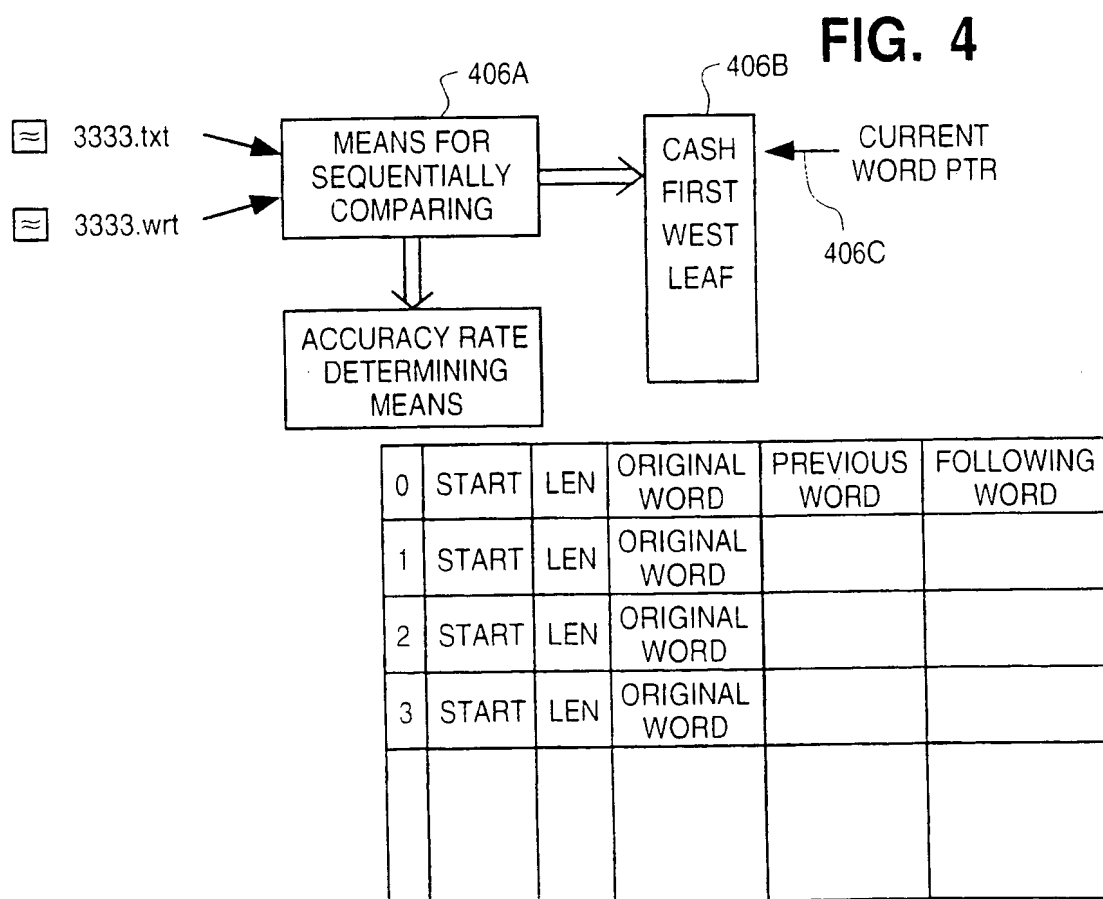
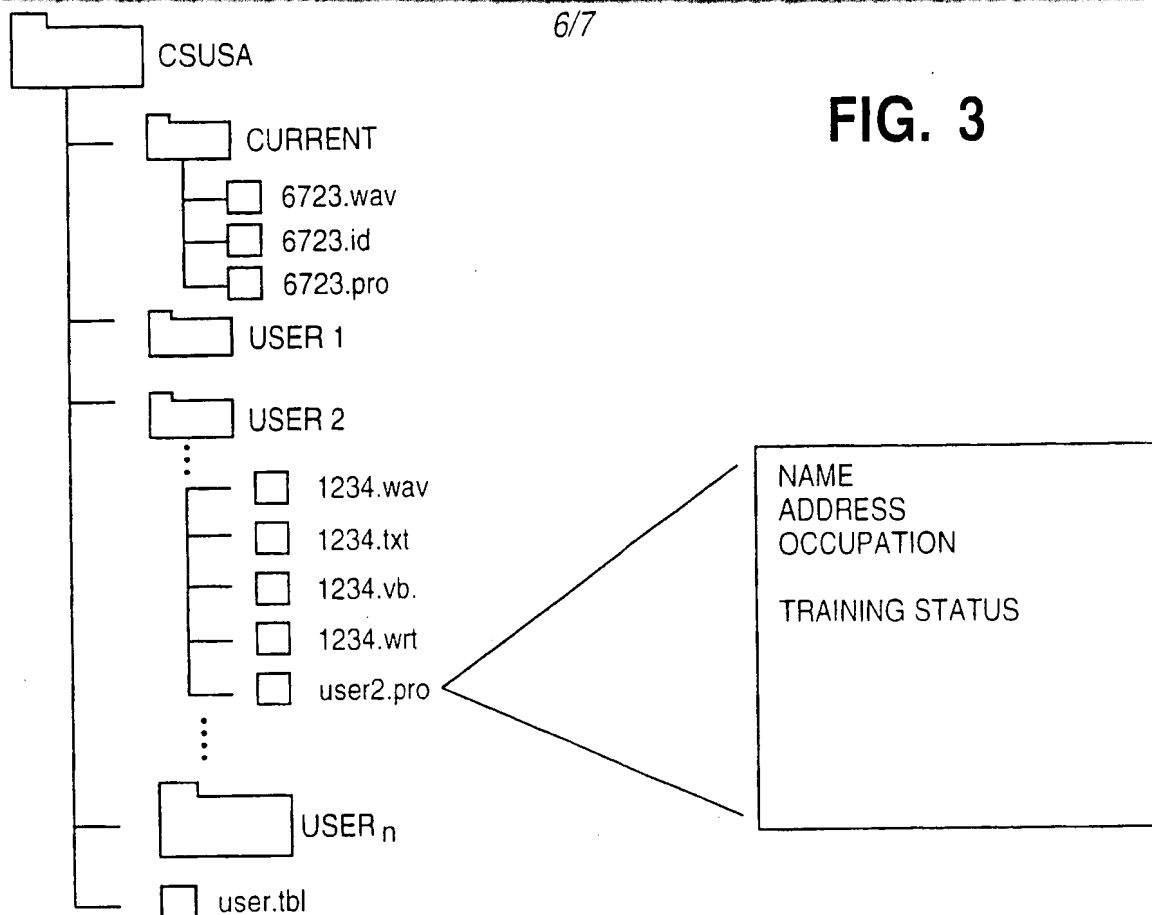
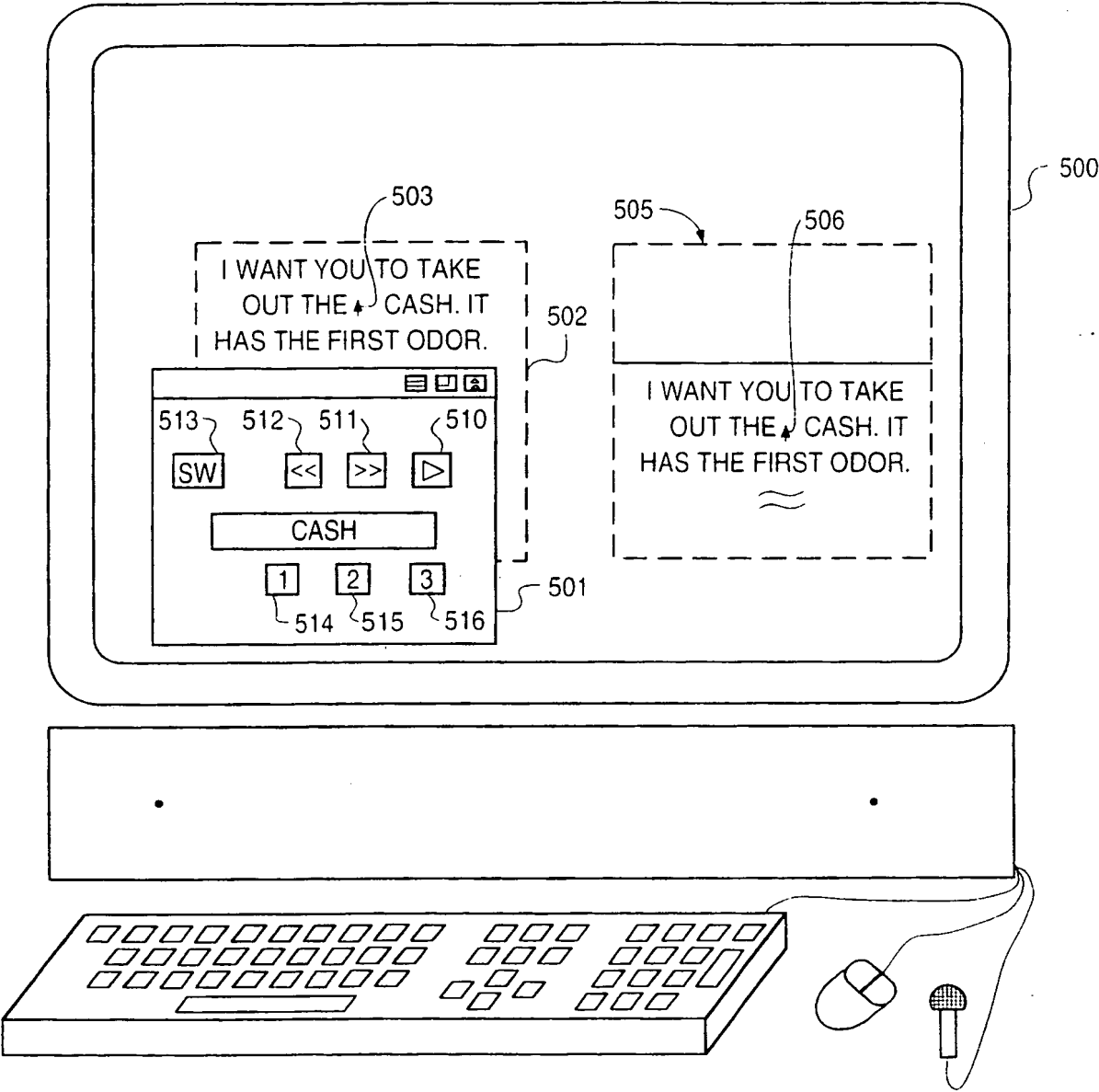


FIG. 5



(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
10 August 2000 (10.08.2000)

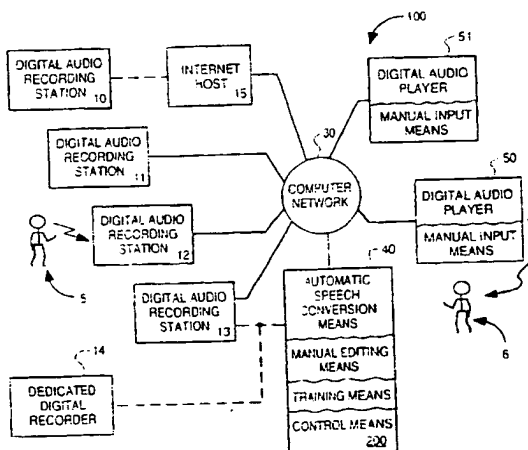
PCT

(10) International Publication Number
WO 00/46787 A3

- (51) International Patent Classification⁷: **G10L 15/06**
- (21) International Application Number: **PCT/US00/02808**
- (22) International Filing Date: **4 February 2000 (04.02.2000)**
- (25) Filing Language: **English**
- (26) Publication Language: **English**
- (30) Priority Data:
60/118,949 **5 February 1999 (05.02.1999)** **US**
- (71) Applicant (for all designated States except US): **CUSTOM SPEECH USA, INC.** [US/US]: Suite B365, 3 North Court Street, Crown Point, IN 46307 (US).
- (72) Inventors: and
- (75) Inventors/Applicants (for US only): **KAHN, Jonathan** [US/US]: 1108 Cheyenne Drive, Crown Point, IN 46307 (US). **QIN, Charles** [—/US]: 23461 North Garden Lane, Lake Zurich, IL 60047 (US). **FLYNN, Thomas, P.** [US/US]: 562 Ridgelawn Road, Crown Point, IN 46307 (US). **TIPPE, Robert, J.** [US/US]: 3813 W. 214th Street, Matteson, IL 60443 (US).
- (74) Agents: **SIGALE, Jordan, A. et al.**: Sonnenschein, Nath & Rosenthal, 8000 Sears Tower, 233 S. Wacker Drive, Chicago, IL 60606-6404 (US).
- (81) Designated States (national): **AE, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CR, CU, CZ, DE, DK, DM, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, TZ, UA, UG, US, UZ, VN, YU, ZA, ZW.**
- (84) Designated States (regional): **ARIPO patent (GH, GM, KE, LS, MW, SD, SL, SZ, TZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).**
- Published:
— With international search report.
- (88) Date of publication of the international search report:
14 December 2000

[Continued on next page]

(54) Title: **SYSTEM AND METHOD FOR AUTOMATING TRANSCRIPTION SERVICES**



(57) Abstract: A system for substantially automating transcription services for multiple users (10, 11, 12) including a manual transcription station (50), speech recognition program (40) and a routing program (200). A uniquely identified voice dictation file is generated from a user and -- based on the training status -- routes the voice dictation file to a manual transcription station and speech recognition program. A human transcriptionist creates transcribed files for each voice dictation file. The speech recognition program creates written text for each dictation file if the training status is training or automated. If the training status of the current user is enrollment or training, a verbatim file is manually established and the speech recognition program is trained with an acoustic model using the verbatim and voice dictation files. The transcribed file is returned to the user if the training status is enrollment or training or written text is returned if the status is automated.

WO 00/46787 A3



For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

INTERNATIONAL SEARCH REPORT

International application No.
PCT/US00/02808

A. CLASSIFICATION OF SUBJECT MATTER

IPC(7) : G10L 15/06

US CL : 704/235, 275

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

U.S. : 704/235, 275

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X -- A	US 5,799,273 A (MITCHELL et al.) 25 August 1998, column 12, lines 15-67 through column 15, lines 1-57. Figures 11-15.	1-16, 20 ----- 17-19, 21-22
Y, P -- A	US 5,875,448 A (BOYS et al.) 23 February 1999, column 15, lines 44-67 and column 16, lines 1-42	1-16, 20 --- 17-19, 21-22
Y -- A	Dragon Dictate for Windows 2.0, User's Guide. British version, First edition, Dragon Systems, Inc., Newton, Massachusetts, pp. 1-230	1-16, 20 ----- 17-19, 21-22

☐ Further documents are listed in the continuation of Box C. ☐ See patent family annex.

* Special categories of cited documents:	*T* later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
A document defining the general state of the art which is not considered to be of particular relevance	*X* document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
B earlier document published on or after the international filing date	*Y* document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
L document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)	*Z* document member of the same patent family
O document referring to an oral disclosure, use, exhibition or other means	
P document published prior to the international filing date but later than the priority date claimed	

Date of the actual completion of the international search 30 JUNE 2000	Date of mailing of the international search report 07 AUG 2000
Name and mailing address of the ISA/US Commissioner of Patents and Trademarks Box PCT Washington, D.C. 20231 Facsimile No. (703) 305-3230	Authorized officer ANGELA ARMSTRONG <i>Joni Hill</i> Telephone No. (703) 305-9000

This Page Blank (uspto)